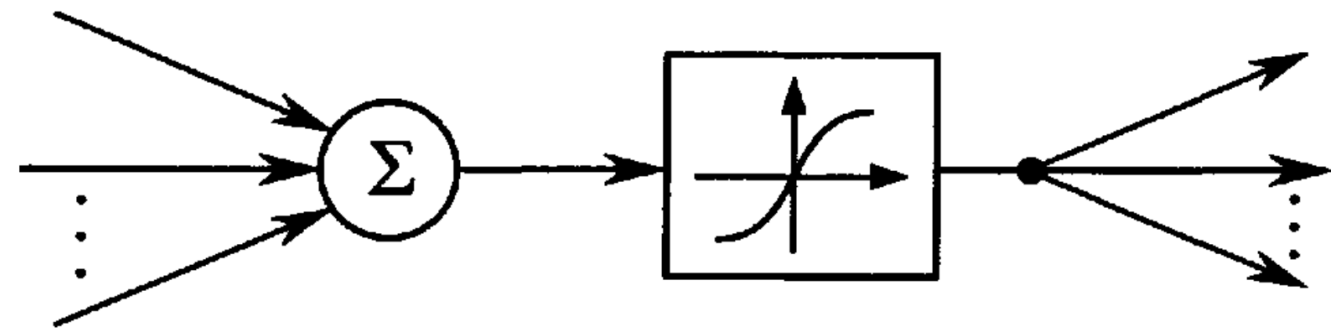


# Unconstrained Optimization and Neural Networks

Jinning Li  
lijinning@sjtu.edu.cn  
jinningli.cn

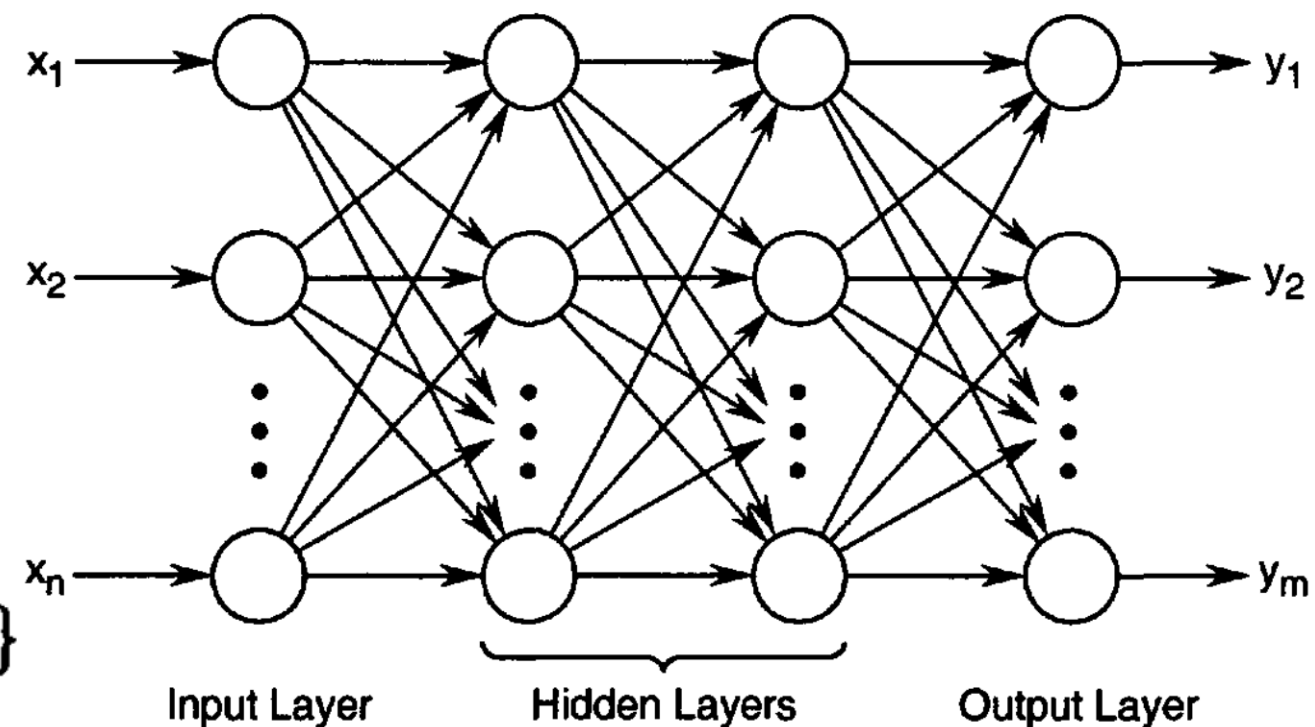
# Introduction

- Single Neuron
- Feedforward Neural network
- A **map** from  $\mathbb{R}^n$  to  $\mathbb{R}^m$ , **n**: number of input, **m**: number of output.



**Figure 13.1** A single neuron

- The information about the mapping is "stored" in the weights over all the neurons
- Training set:  
$$\{(\mathbf{x}_{d,1}, \mathbf{y}_{d,1}), \dots, (\mathbf{x}_{d,p}, \mathbf{y}_{d,p})\}$$
- Function approximators:  
$$\mathbf{y}_{d,i} = \mathbf{F}(\mathbf{x}_{d,i})$$



**Figure 13.3** Structure of a feedforward neural network

# SINGLE-NEURON TRAINING

- Activation function: Identity  
(linear function with unit slope)

$$y = \sum_{i=1}^n w_i x_i = \mathbf{x}^T \mathbf{w}$$

- Training dataset:

$$\mathbf{X}_d = [\mathbf{x}_{d,1} \cdots \mathbf{x}_{d,p}]$$

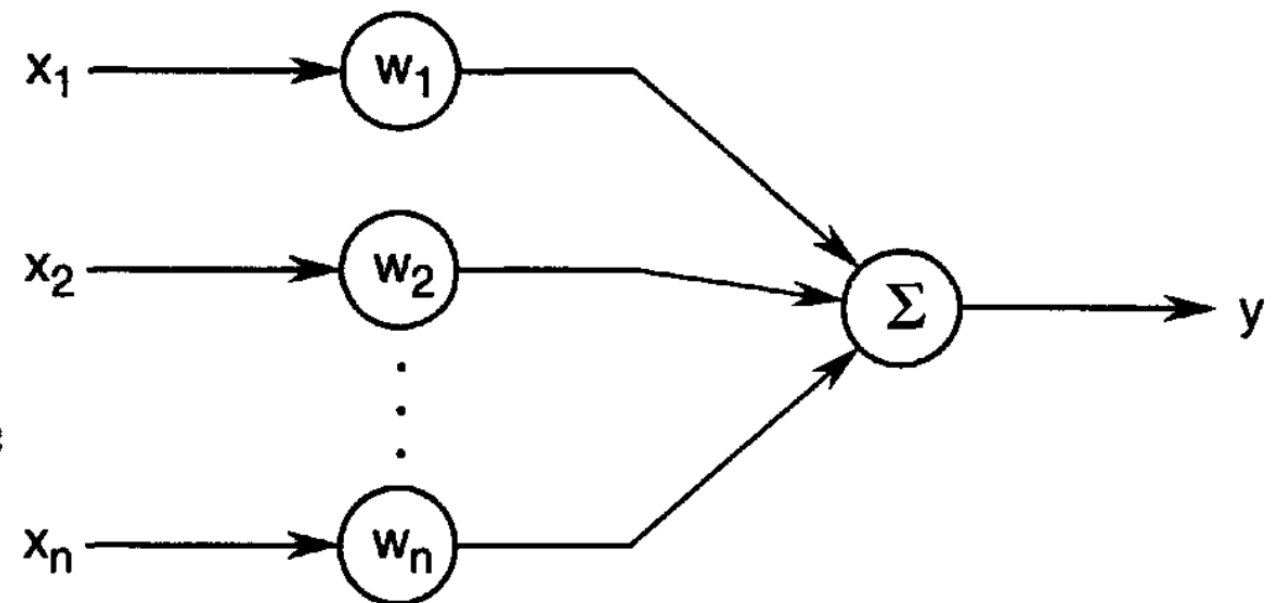
$$\mathbf{y}_d = \begin{bmatrix} y_{d,1} \\ \vdots \\ y_{d,p} \end{bmatrix}.$$

- Objective function:

$$\text{minimize } \frac{1}{2} \sum_{i=1}^p (y_{d,i} - \mathbf{x}_{d,i}^T \mathbf{w})^2$$

- Solution:

- Case 1:  $p \leq n$
- Case 2:  $p > n$



**Figure 13.4** A single linear neuron

# Case 1: $p \leq n$

- Assume  $\text{rank } \mathbf{X}_d^T = p$ , infinite solution to  $\mathbf{y}_d = \mathbf{X}_d^T \mathbf{w}$ .
- Task:  
    minimize  $\|\mathbf{w}\|$   
    subject to  $\mathbf{X}^T \mathbf{w} = \mathbf{y}$
- 其实是个最小二乘问题，上回书说到：

定理 12.1 能够最小化  $\|\mathbf{Ax} - \mathbf{b}\|^2$  的向量  $\mathbf{x}^*$  具有唯一性，可通过求解方程组  $\mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{b}$  得到，即  $\mathbf{x}^* = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{b}$ 。  $\square$

- An iteration method without calculate inverse: Kaczmarz's algorithm:

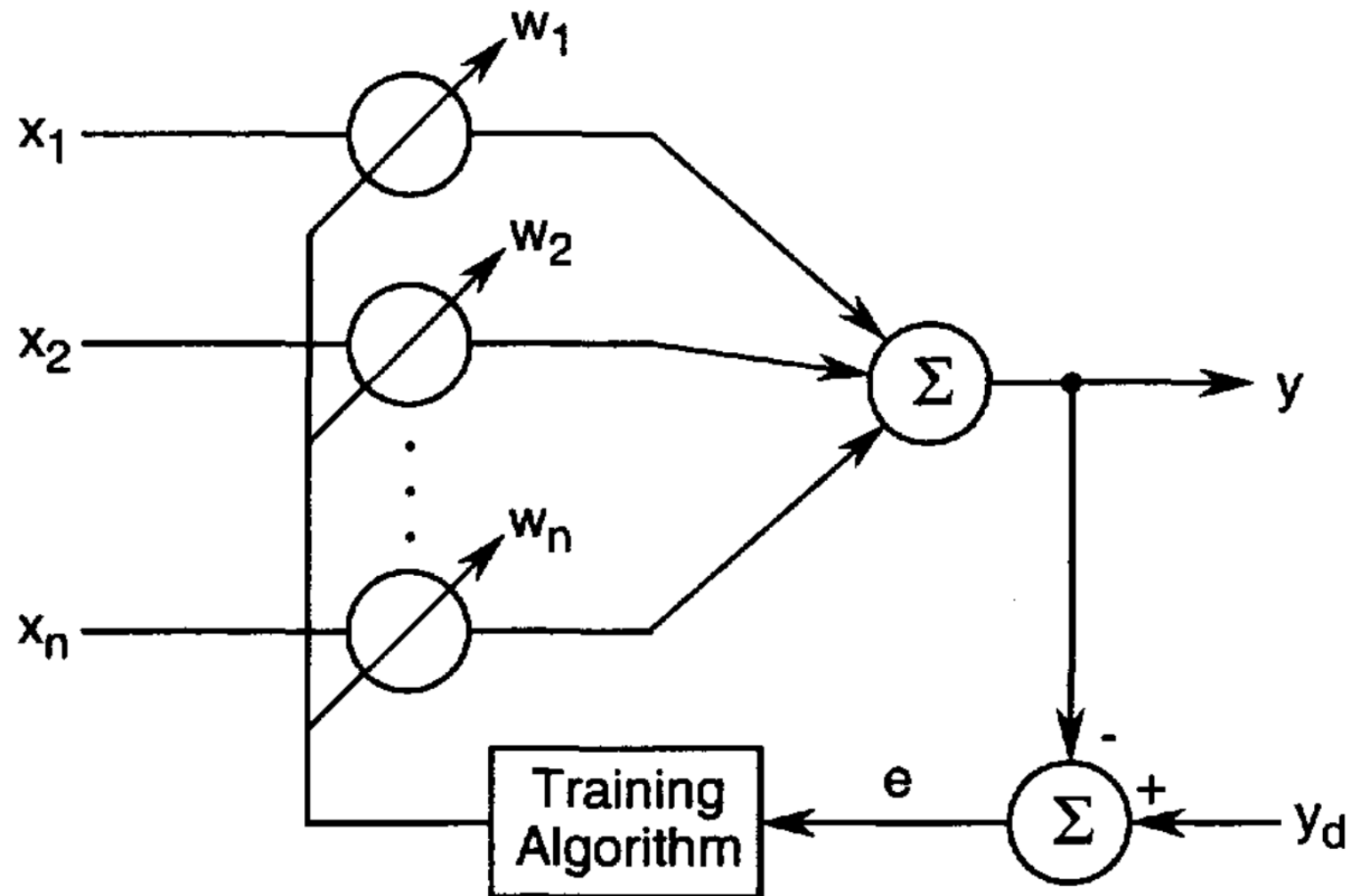
$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \mu \frac{e_k \mathbf{x}_{d,R(k)}}{\|\mathbf{x}_{d,R(k)}\|^2},$$

where  $\mathbf{w}^{(0)} = \mathbf{0}$ , and

$$e_k = y_{d,R(k)} - \mathbf{x}_{d,R(k)}^T \mathbf{w}^{(k)}.$$

# Widrow and Hoff

- Adaline (Adaptive Linear System)



**Figure 13.5** Adaline

# Case 2: $p > n$

- More training points than the number of weights
- $\frac{1}{2} \|\mathbf{y}_d - \mathbf{X}_d^T \mathbf{w}\|^2$  : strictly convex quadratic function of  $\mathbf{w}$ , because  $\mathbf{X}_d \mathbf{X}_d^T$  is a positive definite matrix.
- Unconstrained optimization algorithms: gradient descent

$$\text{minimize } \frac{1}{2} \|\mathbf{y} - \mathbf{X}^T \mathbf{w}\|^2$$

- Objective function:

$$f(\mathbf{w}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}^T \mathbf{w}\|^2$$
$$= \frac{1}{2} (\mathbf{w}^T \mathbf{X}) (\mathbf{X}^T \mathbf{w}) - \mathbf{y}^T \mathbf{X}^T \mathbf{w} + \frac{1}{2} \mathbf{y}^T \mathbf{y}$$

- Derivative:

$$= \frac{1}{2} (\mathbf{w}^T \mathbf{X} \mathbf{X}^T \mathbf{w}) - (\mathbf{X} \mathbf{y})^T \mathbf{w} + \frac{1}{2} \mathbf{y}^T \mathbf{y}$$

$$Df(\mathbf{w}) = \mathbf{X} \mathbf{X}^T \mathbf{w} - \mathbf{X} \mathbf{y}$$

# Case 2: $p > n$

- Iteration formula:  $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \alpha_k \mathbf{X}_d \mathbf{e}^{(k)}$   
where  $\mathbf{e}^{(k)} = \mathbf{y}_d - \mathbf{X}_d^T \mathbf{w}^{(k)}$

- For other activate function:

$$y = f_a \left( \sum_{i=1}^n w_i x_i \right) = f_a (\mathbf{x}^T \mathbf{w})$$

$$\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} + \mu \frac{e_k \mathbf{x}_d}{\|\mathbf{x}_d\|^2}$$

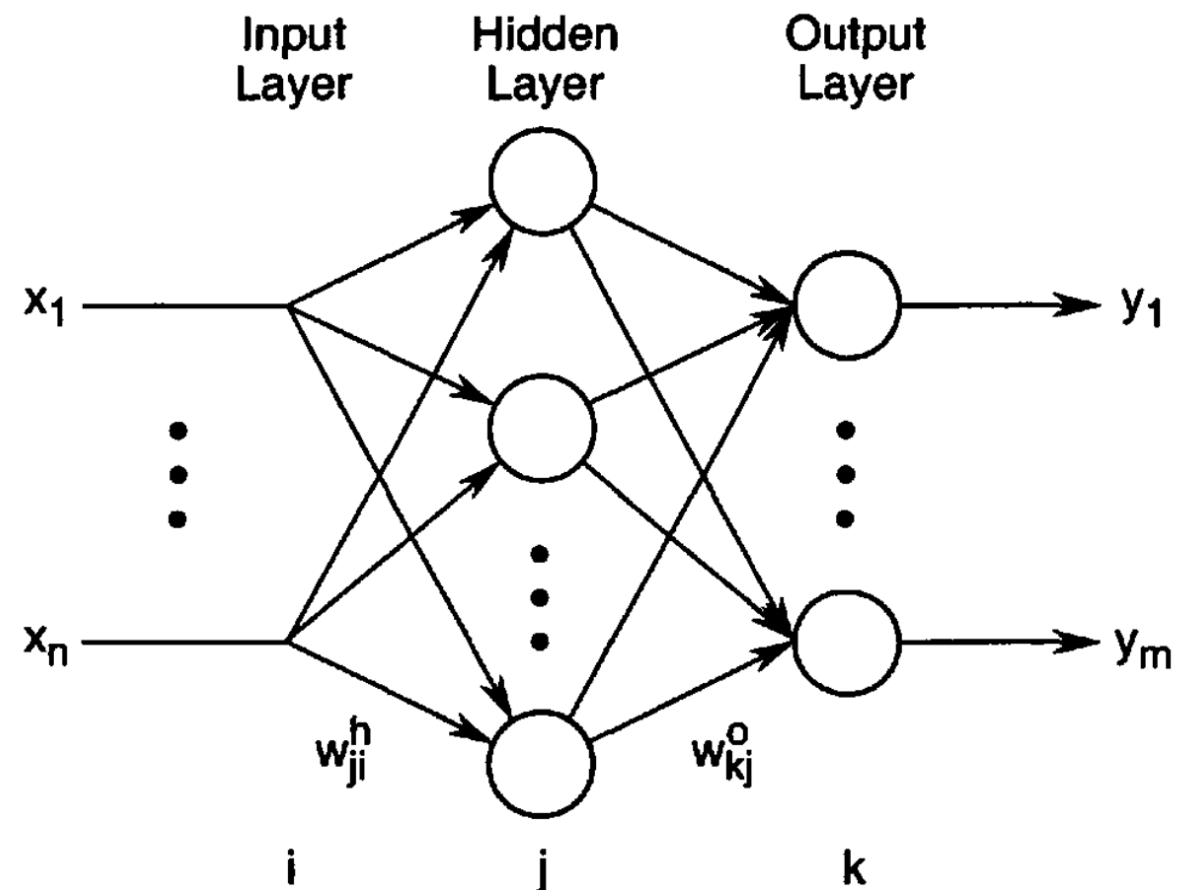
$$e_k = y_d - f_a (\mathbf{x}_d^T \mathbf{w}^{(k)})$$

# BACKPROPAGATION ALGORITHM

- We denote the weights for inputs into the hidden layer by  $w_{ji}^h$
- denote the weights for inputs from the hidden layer into the output layer by  $w_{kj}^o$
- input to the j-th neuron in the hidden layer:  $v_j$
- output of the j-th neuron in the hidden layer:  $z_j$

$$v_j = \sum_{i=1}^n w_{ji}^h x_i,$$

$$z_j = f_j^h \left( \sum_{i=1}^n w_{ji}^h x_i \right)$$



**Figure 13.6** A three-layered neural network



# BACKPROPAGATION ALGORITHM

- The output from the  $s$ th neuron of the output layer is:

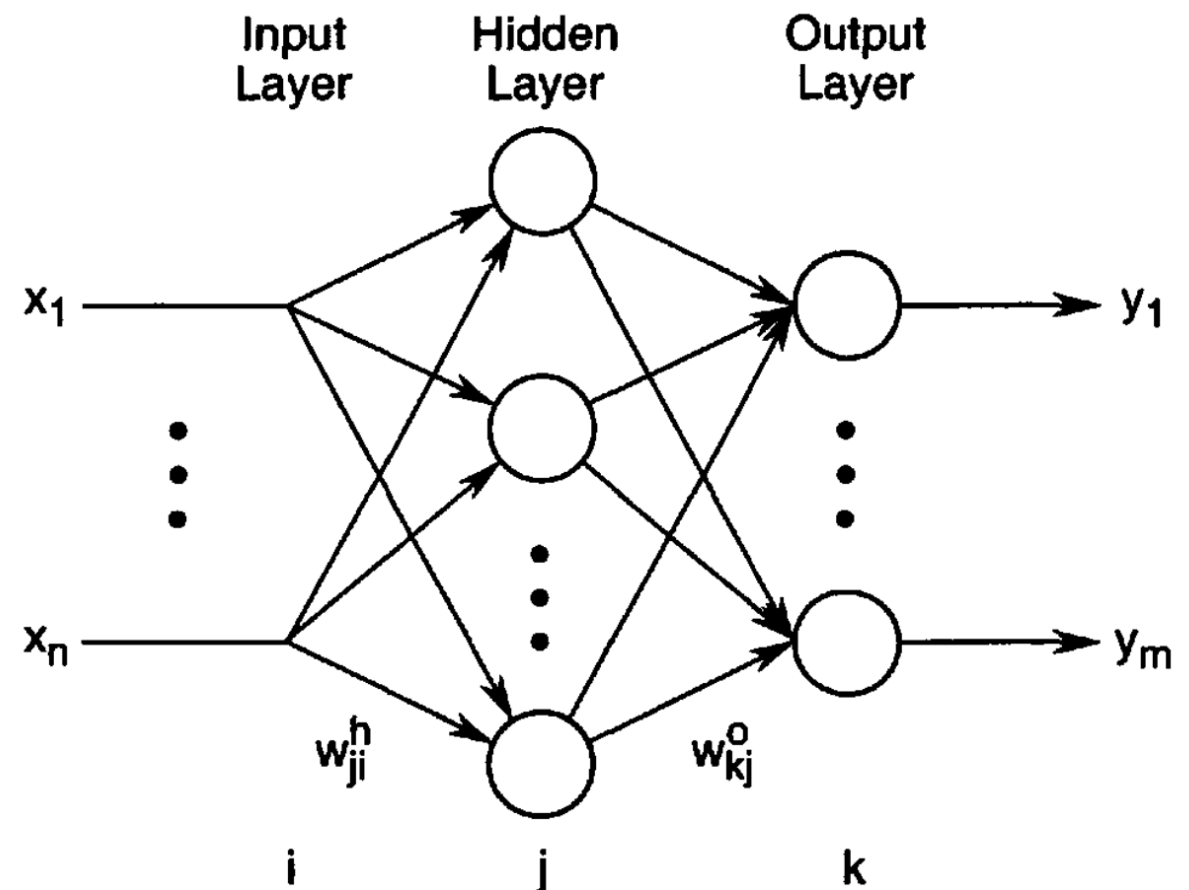
$$y_s = f_s^o \left( \sum_{j=1}^l w_{sj}^o z_j \right)$$

- Therefore, the relationship between the inputs  $\mathbf{x}_i$  and the outputs  $\mathbf{y}_s$  is given by

$$\begin{aligned} y_s &= f_s^o \left( \sum_{j=1}^l w_{sj}^o f_j^h(v_j) \right) \\ &= f_s^o \left( \sum_{j=1}^l w_{sj}^o f_j^h \left( \sum_{i=1}^n w_{ji}^h x_i \right) \right) \\ &= F_s(x_1, \dots, x_n). \end{aligned}$$

- The overall mapping that the neural network implements is therefore given by:

$$\begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} F_1(x_1, \dots, x_n) \\ \vdots \\ F_m(x_1, \dots, x_n) \end{bmatrix}$$



**Figure 13.6** A three-layered neural network

# BACKPROPAGATION ALGORITHM

- Optimization Problem:

$$\text{minimize } \frac{1}{2} \sum_{s=1}^m (y_{ds} - y_s)^2$$

$$y_s = f_s^o \left( \sum_{j=1}^l w_{sj}^o f_j^h \left( \sum_{i=1}^n w_{ji}^h x_i \right) \right)$$

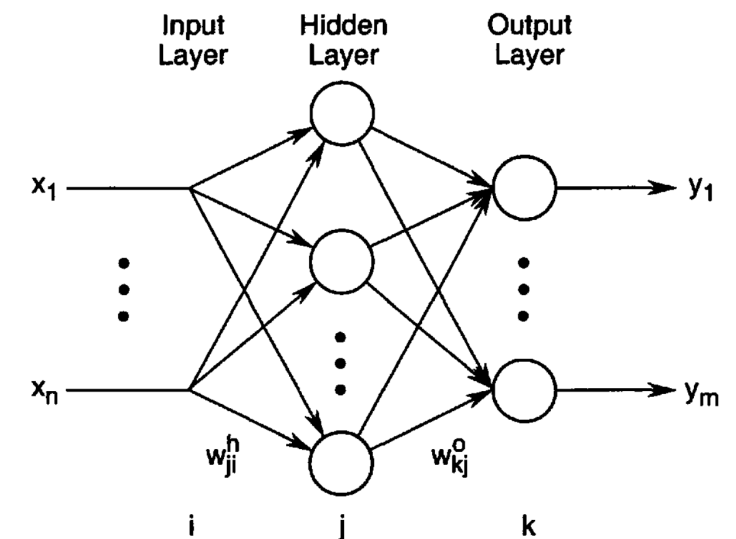


Figure 13.6 A three-layered neural network

- 经过艰难的求导(p226), 得到递推公式:

$$w_{sj}^{o(k+1)} = w_{sj}^{o(k)} + \eta \delta_s^{(k)} z_j^{(k)}$$

$$w_{ji}^{h(k+1)} = w_{ji}^{h(k)} + \eta \left( \sum_{p=1}^m \delta_p^{(k)} w_{pj}^{o(k)} \right) f_j^{h'}(v_j^{(k)}) x_{di},$$

$$v_j^{(k)} = \sum_{i=1}^n w_{ji}^{h(k)} x_{di},$$

$$z_j^{(k)} = f_j^h(v_j^{(k)}),$$

$$y_s^{(k)} = f_s^o \left( \sum_{q=1}^l w_{sq}^{o(k)} z_q^{(k)} \right),$$

$$\delta_s^{(k)} = (y_{ds} - y_s^{(k)}) f_s^{o'} \left( \sum_{q=1}^l w_{sq}^{o(k)} z_q^{(k)} \right)$$



Forward pass: compute the quantities  $v_j^{(k)}$ ,  $z_j^{(k)}$ ,  $y_s^{(k)}$  and  $\delta_s^{(k)}$

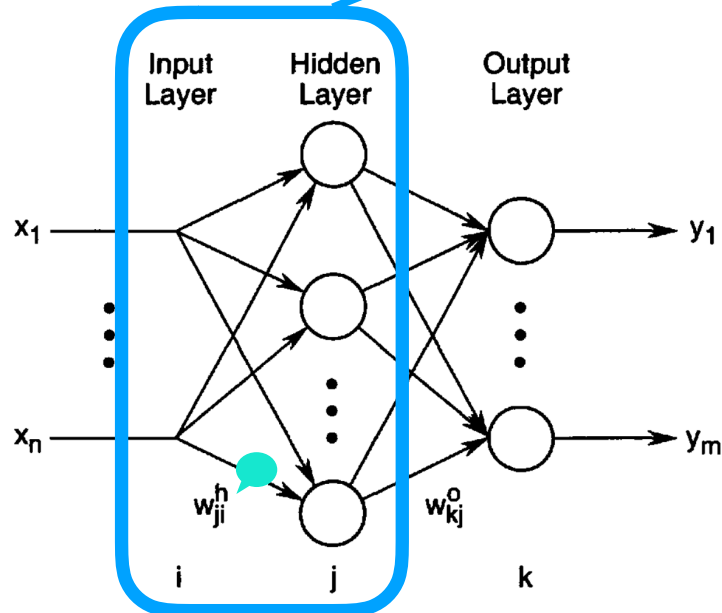


Figure 13.6 A three-layered neural network

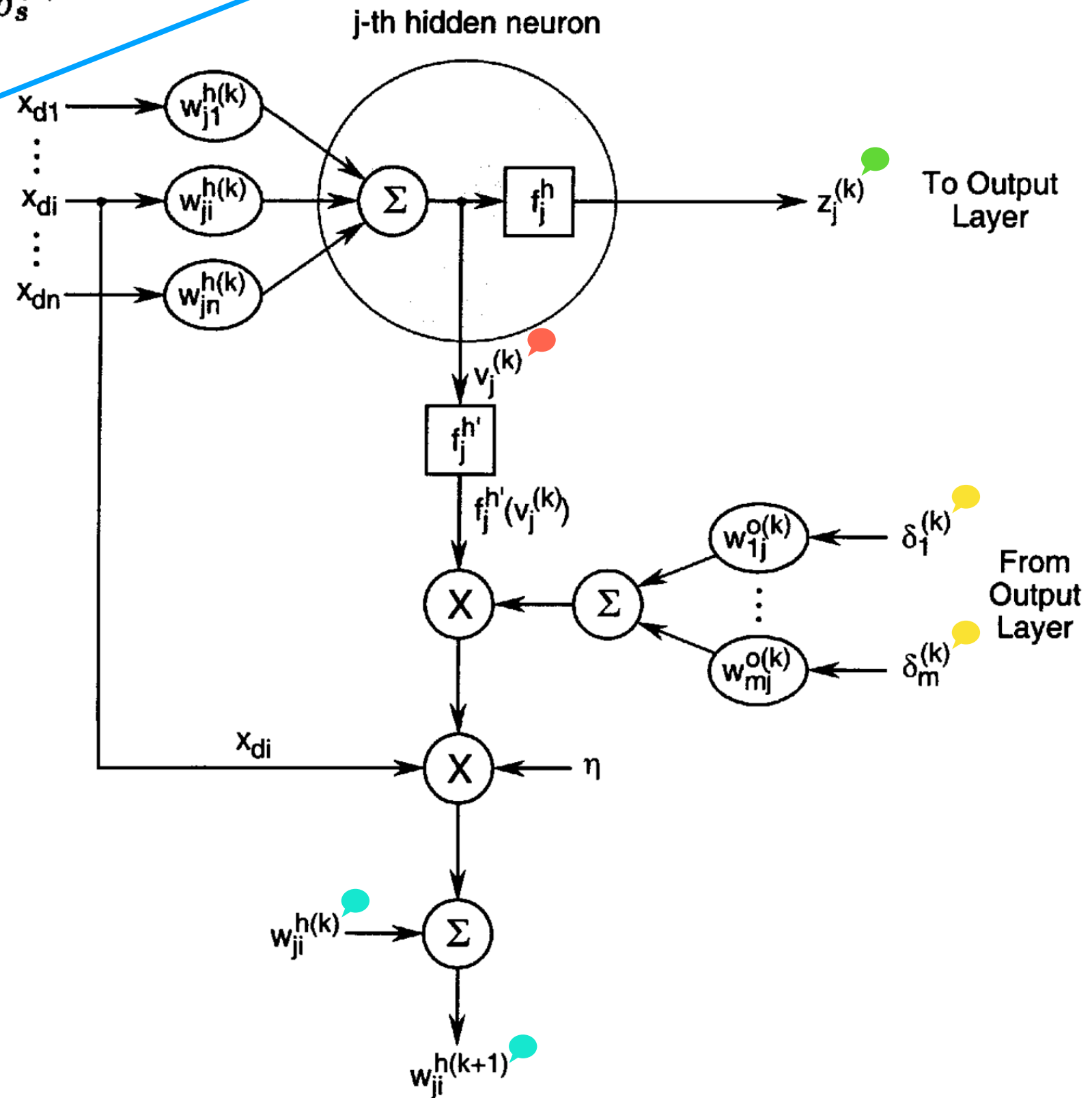


Figure 13.8 Illustration of the update equation for the hidden layer weights

Reverse pass: involves propagating the computed output errors  $\delta_s^{(k)}$  backward

Forward pass: compute the quantities  $v_j^{(k)}$ ,  $z_j^{(k)}$ ,  $y_s^{(k)}$  and  $\delta_s^{(k)}$

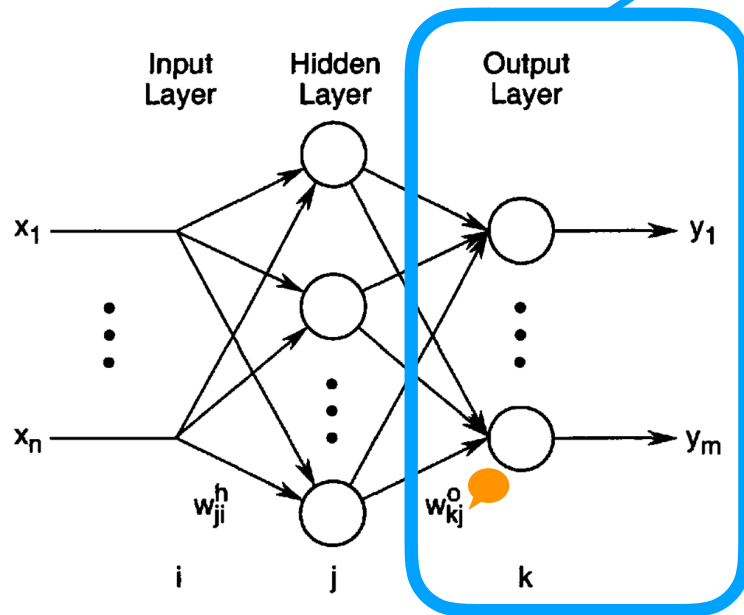


Figure 13.6 A three-layered neural network

Reverse pass: involves propagating the computed output errors  $\delta_s^{(k)}$  backward

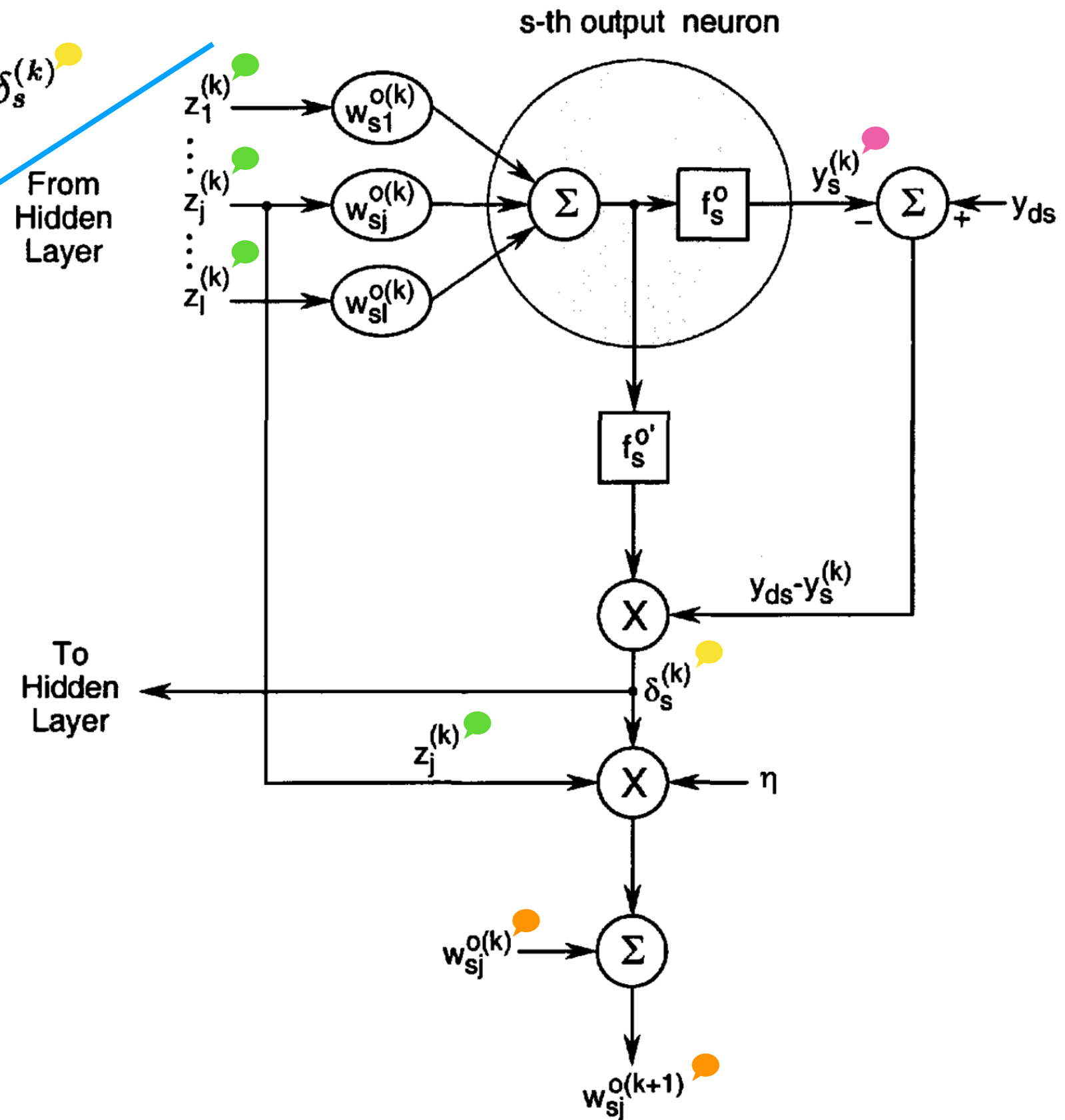
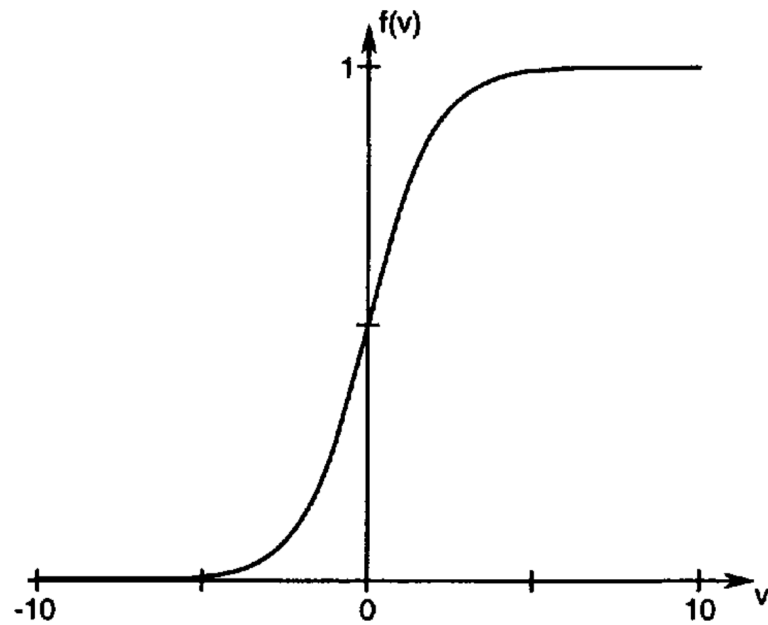


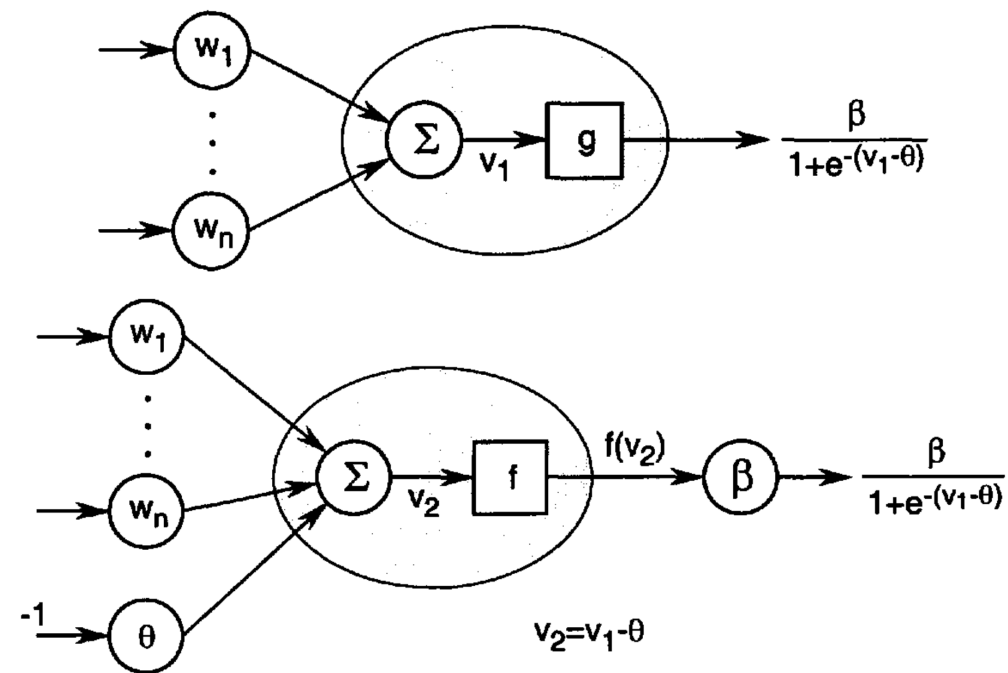
Figure 13.7 Illustration of the update equation for the output layer weights

# Sigmoid Function



**Figure 13.10** The sigmoid function

$$f(v) = \frac{1}{1 + e^{-v}}$$



**Figure 13.11** The above two configurations are equivalent

general version  $g(v) = \frac{\beta}{1 + e^{-(v-\theta)}}$

The parameters  $\beta$  and  $\theta$  represent *scale* and *shift* parameters

parameter  $\theta$  is often interpreted as a threshold

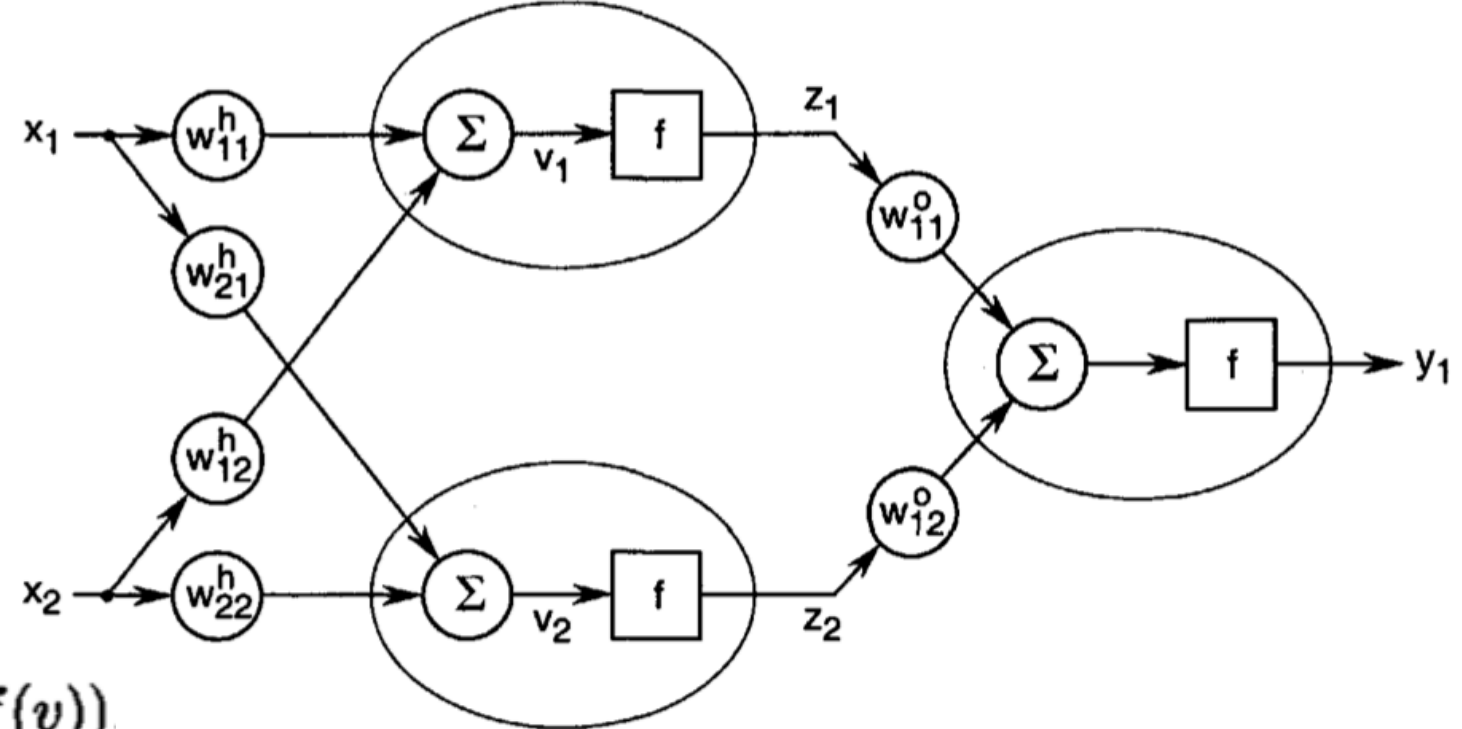
Treating them as additional weights in the backpropagation algorithm. Specifically, we can represent a neuron with activation function  $g$  as one with activation function  $f$  with the addition of two extra weights

# Exercise

- 这是一道填空题。不许看书!



- Active Function:



$$f(v) = 1/(1 + e^{-v}) \quad f'(v) = f(v)(1 - f(v))$$

$$w_{11}^{h(0)} = 0.1 \quad \text{Let } \mathbf{x}_d = [0.2, 0.6]^T$$

$$w_{12}^{h(0)} = 0.3 \quad w_{11}^{o(0)} = 0.4$$

$$w_{21}^{h(0)} = 0.3 \quad w_{12}^{o(0)} = 0.6$$

$$w_{22}^{h(0)} = 0.4 \quad y_d = 0.7 \quad \eta = 10$$

$$v_1^{(0)} = w_{11}^{h(0)} x_{d1} + w_{12}^{h(0)} x_{d2} = \text{(2)一位小数}$$

$$v_2^{(0)} = w_{21}^{h(0)} x_{d1} + w_{22}^{h(0)} x_{d2} = 0.3$$

$$z_1^{(0)} = f(v_1^{(0)}) = \text{(3)表达式} = 0.5498$$

$$z_2^{(0)} = f(v_2^{(0)}) = \frac{1}{1 + e^{-0.3}} = 0.5744$$

$$y_1^{(0)} = f(w_{11}^{o(0)} z_1^{(0)} + w_{12}^{o(0)} z_2^{(0)}) = f(\text{(4)四位小数}) = 0.6375$$

$$\delta_1^{(0)} = (y_d - y_1^{(0)}) y_1^{(0)} (1 - y_1^{(0)}) = 0.01444$$

$$w_{11}^{o(1)} = \text{(5)表达式} = 0.4794$$

$$w_{12}^{o(1)} = w_{12}^{o(0)} + \eta \delta_1^{(0)} z_2^{(0)} = 0.6830$$

$$\delta_1 = (y_d - y_1) f' \left( \sum_{q=1}^2 w_{1q}^o z_q \right)$$

$$= (y_d - y_1) \text{(1)写出表达式}$$

$$= (y_d - y_1) y_1 (1 - y_1)$$

# Thank you!